



THE UNIVERSITY *of* EDINBURGH

## Edinburgh Research Explorer

### A short characterization of relative entropy

**Citation for published version:**

Leinster, T 2019, 'A short characterization of relative entropy', *Journal of mathematical physics*, vol. 60, no. 2, 023302. <https://doi.org/10.1063/1.5026999>

**Digital Object Identifier (DOI):**

[10.1063/1.5026999](https://doi.org/10.1063/1.5026999)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Peer reviewed version

**Published In:**

Journal of mathematical physics

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



# **A short characterization of relative entropy**

Tom Leinster\*

*School of Mathematics, University of Edinburgh, Edinburgh EH9 3FD, Scotland.*

## **Abstract**

We prove characterization theorems for relative entropy (also known as Kullback–Leibler divergence),  $q$ -logarithmic entropy (also known as Tsallis entropy), and  $q$ -logarithmic relative entropy. All three have been characterized axiomatically before, but we show that earlier proofs can be simplified considerably, at the same time relaxing some of the hypotheses.

---

\* <https://www.maths.ed.ac.uk/~tl>; Tom.Leinster@ed.ac.uk

## I. INTRODUCTION

The Shannon entropy of a finite probability distribution  $\mathbf{p} = (p_1, \dots, p_n)$ ,

$$H(\mathbf{p}) = \sum_{i: p_i > 0} p_i \log \frac{1}{p_i},$$

is such an important quantity that many authors have sought short lists of properties that determine  $H$  uniquely. Many such characterization theorems have been found, beginning with one in Shannon's seminal paper of 1948 (Ref. 1, Theorem 2). For instance, Faddeev [2] proved that up to a constant factor,  $H$  is uniquely characterized by symmetry, continuity, and a certain recursivity property.

Accompanying Shannon entropy is the concept of relative entropy, defined as follows. Given probability distributions  $\mathbf{p}$  and  $\mathbf{r}$  on  $n$  elements, the **entropy of  $\mathbf{p}$  relative to  $\mathbf{r}$**  is

$$H(\mathbf{p} \parallel \mathbf{r}) = \sum_{i: p_i > 0} p_i \log \frac{p_i}{r_i} \in [0, \infty].$$

Relative entropy goes by a multitude of names: Kullback–Leibler divergence, directed divergence, discrimination information, relative information, information gain, and so on. In information theory, it measures the wastage when a language whose  $n$  letters have frequencies  $\mathbf{p} = (p_1, \dots, p_n)$  is encoded using a system optimized for a different language with frequencies  $\mathbf{r}$ , instead of the system optimized for the original language. There are other interpretations in other fields, as the plethora of names suggests.

Axiomatic characterizations of relative entropy have also been sought and found. One such theorem is implicit in work of Kannappan and Ng [3]. It states that up to a constant factor, relative entropy is uniquely determined by measurability in each of  $\mathbf{p}$  and  $\mathbf{r}$  separately, invariance under permutations of  $\{1, \dots, n\}$ , the vanishing property  $H(\mathbf{p} \parallel \mathbf{p}) = 0$ , and a certain recursivity equation. (Remark II.7 gives further details.) Their proof was a tour de force of functional equations, involving the solution of the functional equation

$$f(x) + (1-x)g\left(\frac{y}{1-x}\right) = h(y) + (1-y)k\left(\frac{x}{1-y}\right) \quad (1)$$

in four unknown functions, as well as the four-variable functional equation

$$F(x, y) + (1-x)F\left(\frac{u}{1-x}, \frac{v}{1-y}\right) = F(u, v) + (1-u)F\left(\frac{x}{1-u}, \frac{y}{1-v}\right).$$

We give a much simpler proof, at the same time weakening the measurability hypothesis. Our proof involves neither of these equations. Instead, it borrows heavily from a categorical characterization of relative entropy by Baez and Fritz [4], which in turn was inspired by a characterization by Petz [5] of the relative entropy of states of matrix algebras. Our characterization of relative entropy is the first main result, Theorem II.1.

Shannon entropy is just one member (albeit a special one) of a one-parameter family of entropies  $(S_q)_{q \in \mathbb{R}}$ , first investigated by Havrda and Charvát [6] and often misattributed to Tsallis (Remark III.2(ii)). These entropies  $S_q$ , and the accompanying relative entropies, are defined as follows.

For  $q \in \mathbb{R}$ , the  **$q$ -logarithm** is the function  $\ln_q: (0, \infty) \rightarrow \mathbb{R}$  given by

$$\ln_q(x) = \int_1^x t^{-q} dt.$$

The  **$q$ -logarithmic entropy** and  **$q$ -logarithmic relative entropy** are defined by

$$S_q(\mathbf{p}) = \sum_{i: p_i > 0} p_i \ln_q \frac{1}{p_i},$$

$$S_q(\mathbf{p} \parallel \mathbf{r}) = - \sum_{i: p_i > 0} p_i \ln_q \frac{r_i}{p_i},$$

for probability distributions  $\mathbf{p}$  and  $\mathbf{r}$  on  $n$  elements. When  $q = 1$ , these reduce to the ordinary Shannon entropy and relative entropy.

There are several existing theorems characterizing the  $q$ -logarithmic entropy for a given  $q \neq 1$ . Up until now, the simplest appears to have been the 1970 result of Daróczy [7]. We simplify further, weakening the hypotheses and shortening the proof to just a few lines (Theorem III.1). Finally, we use a similar and equally short argument to characterize the  $q$ -logarithmic relative entropies  $S_q(- \parallel -)$  (Theorem IV.1).

It is remarkable that when  $q \neq 1$ , the characterizations of  $q$ -logarithmic entropy and  $q$ -logarithmic relative entropy need no regularity conditions whatsoever (not even measurability), in contrast to the theorems for  $q = 1$ .

The remaining three sections of this paper establish our three theorems in turn, characterizing first relative entropy, then  $q$ -logarithmic entropy, then  $q$ -logarithmic relative entropy.

## II. RELATIVE ENTROPY

For  $n \geq 1$ , write

$$\Delta_n = \left\{ \mathbf{p} = (p_1, \dots, p_n) \in \mathbb{R}^n : p_i \geq 0, \sum p_i = 1 \right\}$$

for the set of probability distributions on  $\{1, \dots, n\}$ , and write

$$A_n = \left\{ (\mathbf{p}, \mathbf{r}) \in \Delta_n \times \Delta_n : p_i = 0 \text{ whenever } r_i = 0 \right\}.$$

Evidently,  $(\mathbf{p}, \mathbf{r}) \in A_n$  if and only if the relative entropy

$$H_n(\mathbf{p} \parallel \mathbf{r}) = \sum_{i: p_i > 0} p_i \log \frac{p_i}{r_i}$$

is finite. (In this section, we add an ‘ $n$ ’ subscript to  $H$  for clarity.) Viewing  $\mathbf{p}$  and  $\mathbf{r}$  as measures on  $\{1, \dots, n\}$ , we have  $(\mathbf{p}, \mathbf{r}) \in A_n$  just when  $\mathbf{p}$  is absolutely continuous with respect to  $\mathbf{r}$ .

We will characterize the sequence of functions

$$H(- \parallel -) = (H_n(- \parallel -) : A_n \rightarrow \mathbb{R})_{n \geq 1}$$

uniquely up to a constant factor. It is easy to check that this sequence has the following four properties, as does any scalar multiple  $cH(- \parallel -)$  with  $c \in \mathbb{R}$ .

**Measurability in the second argument:** For each  $n \geq 1$  and  $\mathbf{p} \in \Delta_n$ , the function

$$\begin{array}{ccc} \{\mathbf{r} \in \Delta_n : (\mathbf{p}, \mathbf{r}) \in A_n\} & \rightarrow & \mathbb{R} \\ \mathbf{r} & \mapsto & H_n(\mathbf{p} \parallel \mathbf{r}) \end{array}$$

is Lebesgue measurable.

**Symmetry:** For each  $n \geq 1$ ,  $(\mathbf{p}, \mathbf{r}) \in A_n$  and permutation  $\sigma$  of  $\{1, \dots, n\}$ ,

$$H_n(\mathbf{p} \parallel \mathbf{r}) = H_n(\mathbf{p}\sigma \parallel \mathbf{r}\sigma), \tag{2}$$

where  $\mathbf{p}\sigma = (p_{\sigma(1)}, \dots, p_{\sigma(n)})$ .

**Vanishing:**  $H_n(\mathbf{p} \parallel \mathbf{p}) = 0$  for all  $n \geq 1$  and  $\mathbf{p} \in \Delta_n$ .

**Chain rule:** To state this, we need some notation. Given  $n, k_1, \dots, k_n \geq 1$  and  $\mathbf{w} \in \Delta_n, \mathbf{p}^1 \in$

$\Delta_{k_1}, \dots, \mathbf{p}^n \in \Delta_{k_n}$ , and writing  $\mathbf{p}^i = (p_1^i, \dots, p_{k_i}^i)$ , define

$$\mathbf{w} \circ (\mathbf{p}^1, \dots, \mathbf{p}^n) = (w_1 p_1^1, \dots, w_1 p_{k_1}^1, \dots, w_n p_1^n, \dots, w_n p_{k_n}^n) \in \Delta_{k_1 + \dots + k_n}.$$

The **chain rule** for relative entropy is that

$$H_{k_1+\dots+k_n}(\mathbf{w} \circ (\mathbf{p}^1, \dots, \mathbf{p}^n) \parallel \tilde{\mathbf{w}} \circ (\tilde{\mathbf{p}}^1, \dots, \tilde{\mathbf{p}}^n)) = H_n(\mathbf{w} \parallel \tilde{\mathbf{w}}) + \sum_{i=1}^n w_i H_{k_i}(\mathbf{p}^i \parallel \tilde{\mathbf{p}}^i) \quad (3)$$

whenever  $(\mathbf{w}, \tilde{\mathbf{w}}) \in A_n$  and  $(\mathbf{p}^i, \tilde{\mathbf{p}}^i) \in A_{k_i}$ . (Under these hypotheses, the pair of distributions on the left-hand side belongs to  $A_{k_1+\dots+k_n}$ .)

**Theorem II.1** *Let  $I(- \parallel -) = (I_n(- \parallel -) : A_n \rightarrow \mathbb{R})_{n \geq 1}$  be a sequence of functions. The following are equivalent:*

- i.  $I(- \parallel -)$  satisfies the four properties above: measurability in the second argument, symmetry, vanishing, and the chain rule;
- ii.  $I(- \parallel -) = cH(- \parallel -)$  for some  $c \in \mathbb{R}$ .

We have just noted that (ii) implies (i). We now embark on the proof of the converse. For the rest of this section, let  $I(- \parallel -) = (I_n(- \parallel -) : A_n \rightarrow \mathbb{R})_{n \geq 1}$  be a sequence of functions satisfying the four conditions. Define a function  $L : (0, 1] \rightarrow \mathbb{R}$  by

$$L(\alpha) = I_2((1, 0) \parallel (\alpha, 1 - \alpha)).$$

The idea is that if  $I(- \parallel -) = H(- \parallel -)$  then  $L = -\log$ . We will show that in any case,  $L$  is a scalar multiple of  $\log$ .

**Lemma II.2** *Let  $(\mathbf{p}, \mathbf{r}) \in A_n$  with  $p_{k+1} = \dots = p_n = 0$ , where  $1 \leq k \leq n$ . Then  $r_1 + \dots + r_k > 0$  and*

$$I_n(\mathbf{p} \parallel \mathbf{r}) = L(r_1 + \dots + r_k) + I_k(\mathbf{p}' \parallel \mathbf{r}'),$$

where

$$\mathbf{p}' = (p_1, \dots, p_k), \quad \mathbf{r}' = \frac{(r_1, \dots, r_k)}{r_1 + \dots + r_k}.$$

**Proof** The case  $k = n$  reduces to the statement that  $L(1) = 0$ , which follows from the vanishing property. Suppose, then, that  $k < n$ .

Since  $\mathbf{p}$  is a probability distribution with  $p_i = 0$  for all  $i > k$ , there is some  $i \leq k$  such that  $p_i > 0$ , and then  $r_i > 0$  as  $(\mathbf{p}, \mathbf{r}) \in A_n$ . Hence  $r_1 + \dots + r_k > 0$ . Let  $\mathbf{r}'' \in \Delta_{n-k}$  be the normalization of  $(r_{k+1}, \dots, r_n)$  if  $r_{k+1} + \dots + r_n > 0$ , or choose  $\mathbf{r}''$  arbitrarily in  $\Delta_{n-k}$  otherwise (which is possible since  $k < n$ ). Then

$$I_n(\mathbf{p} \parallel \mathbf{r}) = I_n((1, 0) \circ (\mathbf{p}', \mathbf{r}'') \parallel (r_1 + \dots + r_k, r_{k+1} + \dots + r_n) \circ (\mathbf{r}', \mathbf{r}'')).$$

The result now follows from the chain rule. □

**Lemma II.3**  $L(\alpha\beta) = L(\alpha) + L(\beta)$  for all  $\alpha, \beta \in (0, 1]$ .

**Proof** We evaluate the real number

$$x := I_3((1, 0, 0) \parallel (\alpha\beta, \alpha(1 - \beta), 1 - \alpha))$$

in two ways. By Lemma II.2 with  $k = 1$  and the vanishing property,

$$x = L(\alpha\beta) + I_1((1) \parallel (1)) = L(\alpha\beta),$$

where  $(1)$  denotes the unique element of  $\Delta_1$ . But also, by Lemma II.2 with  $k = 2$ ,

$$x = L(\alpha) + I_2((1, 0) \parallel (\beta, 1 - \beta)) = L(\alpha) + L(\beta).$$

Comparing the two expressions for  $x$  gives the result.  $\square$

**Lemma II.4** There is some  $c \in \mathbb{R}$  such that  $L(\alpha) = -c \log \alpha$  for all  $\alpha \in (0, 1]$ .

**Proof** Define  $f: [0, \infty) \rightarrow \mathbb{R}$  by  $f(t) = L(e^{-t})$ . By Lemma II.3,  $f$  satisfies Cauchy's functional equation:  $f(t + u) = f(t) + f(u)$  for all  $t, u \in [0, \infty)$ . Also,  $f$  is measurable, since  $L$  is. It is well-known [8] that these conditions force  $f(t) = ct$  for some constant  $c$ , giving  $L(\alpha) = -c \log \alpha$ .  $\square$

Our next lemma is an adaptation of the most ingenious part of Baez and Fritz's argument (Ref. 4, Lemma 4.2).

**Lemma II.5** Let  $(\mathbf{p}, \mathbf{r}) \in A_n$  with  $p_i > 0$  for all  $i \in \{1, \dots, n\}$ . Then  $I_n(\mathbf{p} \parallel \mathbf{r}) = cH_n(\mathbf{p} \parallel \mathbf{r})$ .

**Proof** The hypotheses imply that  $r_i > 0$  for all  $i$ . We can therefore choose some  $\alpha \in (0, 1]$  such that  $r_i - \alpha p_i \geq 0$  for all  $i$ . We will compute the number

$$x := I_{2n}((p_1, \dots, p_n, \underbrace{0, \dots, 0}_n) \parallel (\alpha p_1, \dots, \alpha p_n, r_1 - \alpha p_1, \dots, r_n - \alpha p_n))$$

in two ways. First, by Lemma II.2, Lemma II.4, and the vanishing property,

$$x = L(\alpha) + I_n(\mathbf{p} \parallel \mathbf{p}) = -c \log \alpha.$$

Second, by symmetry, the chain rule, and Lemma II.4,

$$\begin{aligned} x &= I_{2n}((p_1, 0, \dots, p_n, 0) \parallel (\alpha p_1, r_1 - \alpha p_1, \dots, p_n, r_n - \alpha p_n)) \\ &= I_{2n}(\mathbf{p} \circ ((1, 0), \dots, (1, 0)) \parallel \mathbf{r} \circ ((\alpha \frac{p_1}{r_1}, 1 - \alpha \frac{p_1}{r_1}), \dots, (\alpha \frac{p_n}{r_n}, 1 - \alpha \frac{p_n}{r_n}))) \\ &= I_n(\mathbf{p} \parallel \mathbf{r}) + \sum_{i=1}^n p_i L(\alpha \frac{p_i}{r_i}) \\ &= I_n(\mathbf{p} \parallel \mathbf{r}) - c \log \alpha - cH_n(\mathbf{p} \parallel \mathbf{r}). \end{aligned}$$

Comparing the two expressions for  $x$  gives the result.  $\square$

We have now proved that  $I_n(\mathbf{p} \parallel \mathbf{r}) = cH_n(\mathbf{p} \parallel \mathbf{r})$  when  $\mathbf{p}$  has full support. It only remains to prove it for arbitrary  $\mathbf{p}$ .

**Proof of Theorem II.1** Let  $(\mathbf{p}, \mathbf{r}) \in A_n$ . By symmetry, we can assume that  $p_1, \dots, p_k > 0$  and  $p_{k+1} = \dots = p_n = 0$  for some  $k \in \{1, \dots, n\}$ . Writing  $R = r_1 + \dots + r_k$ , we have  $R > 0$  since  $(\mathbf{p}, \mathbf{r}) \in A_n$ , and then

$$I_n(\mathbf{p} \parallel \mathbf{r}) = L(R) + I_k((p_1, \dots, p_k) \parallel \frac{1}{R}(r_1, \dots, r_k))$$

by Lemma II.2. Hence by Lemmas II.4 and II.5,

$$I_n(\mathbf{p} \parallel \mathbf{r}) = -c \log R + cH_k((p_1, \dots, p_k) \parallel \frac{1}{R}(r_1, \dots, r_k)).$$

But by the same argument applied to  $cH$  in place of  $I$  (or by direct calculation), we also have

$$cH_n(\mathbf{p} \parallel \mathbf{r}) = -c \log R + cH_k((p_1, \dots, p_k) \parallel \frac{1}{R}(r_1, \dots, r_k)).$$

The result follows.  $\square$

**Remarks II.6** i. The vanishing axiom cannot be dropped from Theorem II.1. Indeed, the quantity  $\sum_{i: p_i > 0} p_i \log \frac{1}{r_i}$  satisfies the other three axioms but not vanishing.

ii. In the literature on information functions, the chain rule is often replaced by one of two superficially simpler rules. The first is the special case  $k_1 = 2, k_2 = \dots = k_n = 1$ , which is

$$\begin{aligned} H_{n+1}((pw_1, (1-p)w_1, w_2, \dots, w_n) \parallel (\tilde{p}\tilde{w}_1, (1-\tilde{p})\tilde{w}_1, \tilde{w}_2, \dots, \tilde{w}_n)) \\ = H_n(\mathbf{w} \parallel \tilde{\mathbf{w}}) + w_1 H_2((p, 1-p) \parallel (\tilde{p}, 1-\tilde{p})) \end{aligned} \quad (4)$$

$((\mathbf{w}, \tilde{\mathbf{w}}) \in A_n, ((p, 1-p), (\tilde{p}, 1-\tilde{p})) \in A_2)$ . This is known as **recursivity** or **grouping**. The second is the special case  $n = 2$  of the chain rule, which is

$$\begin{aligned} H_{k+\ell}(w\mathbf{p} \oplus (1-w)\mathbf{r} \parallel \tilde{w}\tilde{\mathbf{p}} \oplus (1-\tilde{w})\tilde{\mathbf{r}}) \\ = H_2((w, 1-w) \parallel (\tilde{w}, 1-\tilde{w})) + wH_k(\mathbf{p} \parallel \tilde{\mathbf{p}}) + (1-w)H_\ell(\mathbf{r} \parallel \tilde{\mathbf{r}}), \end{aligned} \quad (5)$$

where

$$w\mathbf{p} \oplus (1-w)\mathbf{r} = (wp_1, \dots, wp_k, (1-w)r_1, \dots, (1-w)r_\ell)$$

and  $((w, 1-w), (\tilde{w}, 1-\tilde{w})) \in A_2, (\mathbf{p}, \tilde{\mathbf{p}}) \in A_k, (\mathbf{r}, \tilde{\mathbf{r}}) \in A_\ell$ . However, straightforward inductions show that in the presence of the symmetry axiom, either one of the special cases (4) or (5) is equivalent to the full chain rule (3). (Similar inductions are carried out in Ref. 9, p. 5–6.) Which to use is, therefore, simply a matter of taste.



**Remark II.7** Here we compare Theorem II.1 with some earlier characterizations of relative entropy. One of the first such theorems was that of Hobson [10], who used stronger hypotheses for the same conclusion. In common with Theorem II.1, he assumed symmetry, vanishing, and the chain rule (in the equivalent form (5)). But he also assumed continuity in both variables (instead of measurability in one) and a monotonicity hypothesis unlike anything in Theorem II.1.

In 1973, Kannappan and Ng [3] proved a result very close to Theorem II.1. They did not *state* that result in Ref. 3, but the closing remarks in another paper by the same authors [11] and the approach of a contemporaneous paper by Kannappan and Rathie [12] strongly suggest the intent. The result was stated explicitly by Csiszár (Ref. 13, Section 2.1), who attributed it to Kannappan and Ng.

There are some slight differences of hypotheses between Kannappan and Ng's theorem and Theorem II.1. They assumed measurability in both variables, whereas we only assumed measurability in the second. (In fact, all we used was that  $I_2((1,0) \parallel -)$  is measurable.) On the other hand, they only needed the vanishing condition for  $(1/2, 1/2)$ , whereas we needed it for all  $\mathbf{p}$ . They used the chain rule in the equivalent form (4). Their proof and ours are entirely different.

### III. $q$ -LOGARITHMIC ENTROPY

Let  $q \in \mathbb{R}$ . The definition of  $q$ -logarithm in the Introduction gives, explicitly,

$$\ln_q(x) = \frac{1}{1-q}(x^{1-q} - 1)$$

for  $x \in (0, \infty)$  and  $q \neq 1$ , while  $\ln_1$  is the natural logarithm  $\log$ . Hence, explicitly, the  $q$ -logarithmic entropy is given by

$$S_q(\mathbf{p}) = \frac{1}{1-q} \left( \sum_{i: p_i > 0} p_i^q - 1 \right)$$

for  $\mathbf{p} \in \Delta_n$  and  $q \neq 1$ , while  $S_1$  is the Shannon entropy  $H$ . We have  $\ln_q(x) \rightarrow \log(x)$  as  $q \rightarrow 1$ , hence also  $S_q(\mathbf{p}) \rightarrow H(\mathbf{p})$  as  $q \rightarrow 1$ .

Fix  $q \in \mathbb{R}$ . The  $q$ -logarithmic entropy satisfies a chain rule

$$S_q(\mathbf{w} \circ (\mathbf{p}^1, \dots, \mathbf{p}^n)) = S_q(\mathbf{w}) + \sum_{i: w_i > 0} w_i^q S_q(\mathbf{p}^i) \quad (6)$$

( $\mathbf{w} \in \Delta_n$ ,  $\mathbf{p}^1 \in \Delta_{k_1}, \dots, \mathbf{p}^n \in \Delta_{k_n}$ ), as is easily checked. In particular, this holds when  $\mathbf{p}^1 = \dots =$

$\mathbf{p}^n = \mathbf{p}$ , say. For  $\mathbf{w} \in \Delta_n$  and  $\mathbf{p} \in \Delta_k$ , write

$$\begin{aligned}\mathbf{w} \otimes \mathbf{p} &= \mathbf{w} \circ (\mathbf{p}, \dots, \mathbf{p}) \\ &= (w_1 p_1, \dots, w_1 p_k, \dots, w_n p_1, \dots, w_n p_k) \in \Delta_{nk}.\end{aligned}$$

In this case, the  $q$ -chain rule (6) gives a  $q$ -**multiplicativity** property:

$$S_q(\mathbf{w} \otimes \mathbf{p}) = S_q(\mathbf{w}) + \left( \sum_{i: w_i > 0} w_i^q \right) S_q(\mathbf{p}) \quad (7)$$

$(n, k \geq 1, \mathbf{w} \in \Delta_n, \mathbf{p} \in \Delta_k)$ .

Note also that  $S_q$  is symmetric in its arguments:

$$S_q(\mathbf{p}) = S_q(\mathbf{p}\sigma) \quad (8)$$

for all  $\mathbf{p} \in \Delta_n$  and permutations  $\sigma$  of  $\{1, \dots, n\}$ .

The left-hand side of equation (7) is symmetric in  $\mathbf{w}$  and  $\mathbf{p}$ , but the right-hand side is not obviously so. This is the key to our second theorem.

**Theorem III.1** *Let  $1 \neq q \in \mathbb{R}$  and let  $I = (I_n: \Delta_n \rightarrow \mathbb{R})_{n \geq 1}$  be a sequence of functions. The following are equivalent:*

- i.  $I$  has the  $q$ -multiplicativity property (7) and the symmetry property (8) (both with  $I$  in place of  $S_q$ );*
- ii.  $I = cS_q$  for some  $c \in \mathbb{R}$ .*

**Proof** By the observations just made, (ii) implies (i). Now assume (i). For all  $\mathbf{w} \in \Delta_n$  and  $\mathbf{p} \in \Delta_k$ , we have  $I_{nk}(\mathbf{w} \otimes \mathbf{p}) = I_{nk}(\mathbf{p} \otimes \mathbf{w})$  by symmetry, so

$$I_n(\mathbf{w}) + \left( \sum_{i: w_i > 0} w_i^q \right) I_k(\mathbf{p}) = I_k(\mathbf{p}) + \left( \sum_{i: p_i > 0} p_i^q \right) I_n(\mathbf{w}),$$

or equivalently

$$\left( \sum_{i: w_i > 0} w_i^q - 1 \right) I_k(\mathbf{p}) = \left( \sum_{i: p_i > 0} p_i^q - 1 \right) I_n(\mathbf{w}).$$

Take  $\mathbf{w} = (1/2, 1/2)$ : then for all  $\mathbf{p} \in \Delta_k$ ,

$$(2^{1-q} - 1) I_k(\mathbf{p}) = \left( \sum_{i: p_i > 0} p_i^q - 1 \right) I_2(1/2, 1/2).$$

Since  $q \neq 1$ , we can define  $c = \frac{1-q}{2^{1-q}-1} \cdot I_2(1/2, 1/2)$ , and then  $I = cS_q$ . □

**Remarks III.2** i. The  $q$ -logarithms were used in Hardy, Littlewood and Pólya's classic book on inequalities, first published in 1934 (Ref. 14, proof of Theorem 84). They have been an explicit object of study since at least a 1964 paper of Box and Cox in statistics (Ref. 15, Section 3). The name ' $q$ -logarithm' appears to have been introduced by Umarov, Tsallis and Steinberg [16] in 2008, working in statistical mechanics.

ii. The  $q$ -logarithmic entropies have been discovered and rediscovered repeatedly. They seem to have first appeared in a 1967 paper on information and classification by Havrda and Charvát [6], who used a form adapted to base 2 logarithms. They were rediscovered in 1970 by Daróczy [7]. The base  $e$  version  $S_q$  appeared in a 1982 article of Patil and Taillie (Ref. 17, Section 3.2), where it was studied as an index of biodiversity.

In physics, meanwhile, the  $q$ -logarithmic entropies appeared in a 1971 article of Lindhard and Nielsen [18] (according to Csiszár: Ref. 13, Section 2.4), and in a 1978 survey by Wehrl (Ref. 19, p. 247). Finally, they were rediscovered again in a 1988 paper on statistical physics by Tsallis [20].

Despite the twenty years of active life that the  $q$ -logarithmic entropies had already enjoyed, it is after Tsallis that they are most commonly named. The term ' $q$ -logarithmic entropy' is new, but has the benefits of being descriptive and of not perpetuating a misattribution.

iii. As in Remark II.6(ii), a simple inductive argument shows that the  $q$ -chain rule of equation (6) follows from the special case

$$S_q(pw_1, (1-p)w_1, w_2, \dots, w_n) = S_q(\mathbf{w}) + w_1^q S_q(p, 1-p) \quad (9)$$

$$(p \in [0, 1], n \geq 1, \mathbf{w} \in \Delta_n).$$

iv. A characterization of the  $q$ -logarithmic entropies similar to Theorem III.1 was published by Daróczy in 1970 [7]. He assumed the full  $q$ -chain rule for  $I(\mathbf{w} \circ (\mathbf{p}^1, \dots, \mathbf{p}^n))$  (in the equivalent form (9)), rather than just the special case  $\mathbf{p}^1 = \dots = \mathbf{p}^n$  used here. However, where we assumed that  $I_n: \Delta_n \rightarrow \mathbb{R}$  is symmetric for all  $n \geq 2$ , Daróczy only assumed it for  $n = 3$ . The two proofs are very different; the main step in Daróczy's was the solution of the functional equation (1) in the case  $f = g = h = k$ .

Other characterizations of  $S_q$  have been proved, using stronger hypotheses than Theorem III.1 to obtain the same conclusion (such as the theorem in Section 2 of Ref. 21, and Theorem V.2 of Ref. 22).

#### IV. $q$ -LOGARITHMIC RELATIVE ENTROPY

For  $q \neq 1$ , the  $q$ -logarithmic relative entropy  $S_q: A_n \rightarrow \mathbb{R}$ , defined in the Introduction, is given explicitly by

$$S_q(\mathbf{p} \parallel \mathbf{r}) = \frac{1}{q-1} \left( \sum_{i: p_i > 0} p_i^q r_i^{1-q} - 1 \right)$$

for  $(\mathbf{p}, \mathbf{r}) \in A_n$ . In the case  $q = 1$ , it reduces to the ordinary relative entropy  $H(\mathbf{p} \parallel \mathbf{r})$ . As in that case, restricting the arguments to lie in  $A_n$  guarantees that  $S_q(- \parallel -)$  takes only finite values.

Our third and final theorem is a characterization of  $q$ -logarithmic relative entropy, very similar to the characterization of  $q$ -logarithmic entropy itself.

We begin by noting two properties of  $q$ -logarithmic relative entropy. First, there is an easily-checked chain rule:

$$S_q(\mathbf{w} \circ (\mathbf{p}^1, \dots, \mathbf{p}^n) \parallel \tilde{\mathbf{w}} \circ (\tilde{\mathbf{p}}^1, \dots, \tilde{\mathbf{p}}^n)) = S_q(\mathbf{w} \parallel \tilde{\mathbf{w}}) + \sum_{i: w_i > 0} w_i^q \tilde{w}_i^{1-q} S_q(\mathbf{p}^i \parallel \tilde{\mathbf{p}}^i)$$

$((\mathbf{w}, \tilde{\mathbf{w}}) \in A_n, (\mathbf{p}^i, \tilde{\mathbf{p}}^i) \in A_{k_i})$ . This specializes to a  $q$ -multiplicativity formula

$$S_q(\mathbf{w} \otimes \mathbf{p} \parallel \tilde{\mathbf{w}} \otimes \tilde{\mathbf{p}}) = S_q(\mathbf{w} \parallel \tilde{\mathbf{w}}) + \left( \sum_{i: w_i > 0} w_i^q \tilde{w}_i^{1-q} \right) S_q(\mathbf{p} \parallel \tilde{\mathbf{p}}) \quad (10)$$

$((\mathbf{w}, \tilde{\mathbf{w}}) \in A_n, (\mathbf{p}, \tilde{\mathbf{p}}) \in A_k)$ . Second,  $q$ -logarithmic relative entropy has the same symmetry property as ordinary relative entropy:

$$S_q(\mathbf{p} \parallel \mathbf{r}) = S_q(\mathbf{p}\sigma \parallel \mathbf{r}\sigma) \quad (11)$$

for all  $n \geq 1$ ,  $(\mathbf{p}, \mathbf{r}) \in A_n$ , and permutations  $\sigma$  of  $\{1, \dots, n\}$ .

**Theorem IV.1** *Let  $1 \neq q \in \mathbb{R}$  and let  $I(- \parallel -) = (I_n(- \parallel -): A_n \rightarrow \mathbb{R})_{n \geq 1}$  be a sequence of functions. The following are equivalent:*

- i.  $I(- \parallel -)$  has the  $q$ -multiplicativity property (10) and the symmetry property (11) (both with  $I$  in place of  $S_q$ );
- ii.  $I(- \parallel -) = cS_q(- \parallel -)$  for some  $c \in \mathbb{R}$ .

**Proof** It is trivial that (ii) implies (i). Now assume (i). By symmetry,

$$I_{nk}(\mathbf{w} \otimes \mathbf{p} \parallel \tilde{\mathbf{w}} \otimes \tilde{\mathbf{p}}) = I_{nk}(\mathbf{p} \otimes \mathbf{w} \parallel \tilde{\mathbf{p}} \otimes \tilde{\mathbf{w}})$$

for all  $n, k \geq 1$ ,  $(\mathbf{w}, \tilde{\mathbf{w}}) \in A_n$ , and  $(\mathbf{p}, \tilde{\mathbf{p}}) \in A_k$ . So by  $q$ -multiplicativity,

$$I_n(\mathbf{w} \parallel \tilde{\mathbf{w}}) + \left( \sum_{i: w_i > 0} w_i^q \tilde{w}_i^{1-q} \right) I_k(\mathbf{p} \parallel \tilde{\mathbf{p}}) = I_k(\mathbf{p} \parallel \tilde{\mathbf{p}}) + \left( \sum_{i: p_i > 0} p_i^q \tilde{p}_i^{1-q} \right) I_n(\mathbf{w} \parallel \tilde{\mathbf{w}}),$$

or equivalently,

$$\left( \sum_{i: w_i > 0} w_i^q \tilde{w}_i^{1-q} - 1 \right) I_k(\mathbf{p} \parallel \tilde{\mathbf{p}}) = \left( \sum_{i: p_i > 0} p_i^q \tilde{p}_i^{1-q} - 1 \right) I_n(\mathbf{w} \parallel \tilde{\mathbf{w}}).$$

Take  $\mathbf{w} = (1, 0)$  and  $\tilde{\mathbf{w}} = (1/2, 1/2)$ : then

$$(2^{q-1} - 1) I_k(\mathbf{p} \parallel \tilde{\mathbf{p}}) = \left( \sum_{i: p_i > 0} p_i^q \tilde{p}_i^{1-q} - 1 \right) I_2((1, 0) \parallel (1/2, 1/2))$$

for all  $(\mathbf{p}, \tilde{\mathbf{p}}) \in A_k$ . But  $q \neq 1$ , so we can define

$$c = \frac{1-q}{2^{q-1}-1} \cdot I_2((1, 0) \parallel (1/2, 1/2)),$$

and then  $I(- \parallel -) = c S_q(- \parallel -)$ . □

**Remarks IV.2** i. The definition of  $q$ -logarithmic relative entropy was given in 1972 by Rathie and Kannappan [23] (who used a version adapted to base 2 logarithms). The base  $e$  version used here was studied by Cressie and Read in a 1984 paper in statistics (Ref. 24, Section 5). It was rediscovered in physics in 1998, by Shiino [25] and Tsallis [26] independently.

ii. Other characterization theorems for  $q$ -logarithmic relative entropy have been proved. For example, Furuichi (Ref. 22, Section IV) obtained the same conclusion, but also assumed continuity and essentially the full chain rule (that is, an equivalent special case, as in Remarks II.6(ii) and III.2(iii)).

*Acknowledgements* I thank John Baez and Tobias Fritz for their comments. This work was partially supported by a BBSRC FLIP award (BB/P004210/1).

- 
- [1] C. E. Shannon, Bell System Technical Journal **27**, 379 (1948).
  - [2] D. K. Faddeev, Uspekhi Matematicheskikh Nauk **11**, 227 (1956).
  - [3] P. Kannappan and C. T. Ng, Proceedings of the American Mathematical Society **38**, 303 (1973).
  - [4] J. Baez and T. Fritz, Theory and Applications of Categories **29**, 421 (2014).

- [5] D. Petz, *Acta Mathematica Hungarica* **59**, 449 (1992).
- [6] J. Havrda and F. Charvát, *Kybernetika* **3**, 30 (1967).
- [7] Z. Daróczy, *Information and Control* **16**, 36 (1970).
- [8] S. Banach, *Fundamenta Mathematicae* **1**, 123 (1920).
- [9] A. Feinstein, *Foundations of Information Theory* (McGraw–Hill, New York, 1958).
- [10] A. Hobson, *Journal of Statistical Physics* **1**, 383 (1969).
- [11] P. Kannappan and C. T. Ng, *Pacific Journal of Mathematics* **54**, 157 (1974).
- [12] P. Kannappan and P. N. Rathie, *Information and Control* **22**, 163 (1973).
- [13] I. Csiszár, *Entropy* **10**, 261 (2008).
- [14] G. Hardy, J. E. Littlewood, and G. Pólya, *Inequalities*, 2nd ed. (Cambridge University Press, Cambridge, 1952).
- [15] G. E. P. Box and D. R. Cox, *Journal of the Royal Statistical Society. Series B (Methodological)* **26**, 211 (1964).
- [16] S. Umarov, C. Tsallis, and S. Steinberg, *Milan Journal of Mathematics* **76**, 307 (2008).
- [17] G. P. Patil and C. Taillie, *Journal of the American Statistical Association* **77**, 548 (1982).
- [18] J. Lindhard and V. Nielsen, *Matematisk-Fysiske Meddelelser: Kongelige Danske Videnskabernes Selskab* **38**, 1 (1971).
- [19] A. Wehrl, *Reviews of Modern Physics* **50**, 221 (1978).
- [20] C. Tsallis, *Journal of Statistical Physics* **52**, 479 (1988).
- [21] H. Suyari, *Journal of Physics A: Mathematical and General* **35**, 10731 (2002).
- [22] S. Furuichi, *IEEE Transactions on Information Theory* **51**, 3638 (2005).
- [23] P. N. Rathie and P. Kannappan, *Information and Control* **20**, 38 (1972).
- [24] N. Cressie and T. R. C. Read, *Journal of the Royal Statistical Society. Series B (Methodological)* **46**, 440 (1984).
- [25] M. Shiino, *Journal of the Physical Society of Japan* **67**, 3658 (1998).
- [26] C. Tsallis, *Physical Review E* **58**, 1442 (1998).